

Semantic Web構築に向けたメタデータ自動作成に関する研究

著者	岡野 拓生, 鈴木 育男, 久保 洋
雑誌名	サテライト・ベンチャー・ビジネス・ラボラトリー 年報
巻	7
ページ	62-63
発行年	2005
URL	http://hdl.handle.net/10258/312

Semantic Web構築に向けたメタデータ自動作成に関する研究

岡野拓生 (B4)[†], 鈴木育男[‡], 久保洋[†]

[†]室蘭工業大学 情報工学科

[‡]室蘭工業大学 SVBL

1. はじめに

ネットワーク空間での情報分散化を重視したハイパーテキストシステムの一つである Web は、記述の容易さやコンピュータの普及、インターネット人口の増加などから爆発的に普及し、誰もがアクセスできる巨大な情報共有空間となっている。毎日の天気や最新のニュースなど、Web から提供される情報は、我々の生活にはなくてはならない社会的基盤の一つとなってきたと考えられる。

しかし、Web 上に構築される情報空間が巨大になるにつれて、ハイパーリンクを辿る事によって情報を獲得するだけでは、目的とする情報を獲得するまでに多くの労力を必要とし、効率が悪い。そのため、この情報検索の労力を減少させるために、これまでに Yahoo や Google といった情報検索サービスが提供されてきた。このような情報検索サービスでは、検索のための入力文字の選択が検索効率に大きな影響を与えている。つまり、検索サービスを使用するユーザが、適切な検索文字を入力しないと目的の情報が入手できないことになる。これは、現在使用されている Web 上の情報記述言語である HTML が、文字サイズや段組などの装飾情報しかもっていないことが原因であると考えられる。Web 上の情報に、その意味や概念等の情報を付加できれば、その概念情報を基によりよい情報提供が可能となる。

このような問題を解決するための一つの方法として、「セマンティック Web」が提唱されている。このセマンティック Web では、これまでの HTML 文書処理するばかりではなく、その情報の意味内容を表現するメタデータを付加し、このメタデータに対して様々な情報処理を行うことで Web 上の情報を効果的に利用しようという試みである。またメタデータ自体は、HTML 文書以外の情報にも関連づけることが可能であるため、画像・動画・音楽といった文字以外の情報の活用も可能になる。

セマンティック Web では、メタデータを基本としているので、メタデータ情報が普及しなければセマンティック Web の実現はできない。セマンティック Web の実現に向けたもっとも根本的な問題として、「誰がメタデータを書くのか」ということがあげられる。これは、これから新しく発信される情報に対してメタデータを付加するということというよりは、むしろ、現在 Web 上に存在している膨大な情報にメタデータを付加して再構築するかという問題のほうが大きいように感じる。

そこで、本研究ではセマンティック Web の実現に向けて、メタデータを容易に付加できるシステム環境を構築することを目的とする。具体的には、現存する HTML 文書からメタデータを必要とする部分を抜き出し、メタデータ用のタグ情報を自動的に追加するものとする。

2. セマンティック Web 技術と関連研究

セマンティック Web では、前述したように Web ページ (HTML 文書) にそのページの内容や関連した情報について、機械的な処理が可能なメタデータを付加することにより、情報の質をコントロールすることを目指している。さらに、付加したメタデータの意味や相互関係をオントロジーによって関連づけることで、メタデータ間の概念関係を表現可能となり、それによって表層語の違い (日本語と英語の違いなど) を補うことができる。

セマンティック Web において、メタデータの記述言語には、RDF(Resource Description Framework)が標準化されている。RDF の意味モデルは、主語 (subject)、述語 (predicate)、目的語 (object) の 3 つの要素により構成されており、この主語と目的語を述語で結合する有向グラフとして表現される。例として、『お好み焼き』の『材料』が『キャベツ』であることを RDF 意味グラフでは、図 1 の有向グラフで表現される。

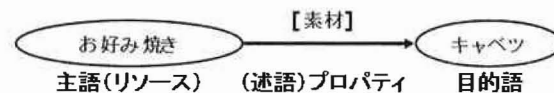


図 1: RDF モデル

このように、RDF を利用すれば、Web 上に意味ネットワークを構築できるのであるが、RDF の作成は容易な作業ではない。現在、メタデータ自動作成機能としては、南野ら[1]の「何でも RSS」, 「MyRSS」[2]などにおいて、RDF の機能を絞った RSS の記述に関して研究が行われている。

3. 提案手法

本研究では、HTML から RDF 記述を自動的に行えるシステムの構築を試みる。構築システムの処理作業は以下のようなになる。

1. HTML 文書を形態素解析によりメタデータを付加する要素を抽出
2. 抽出した要素に専用のタグを付加
3. タグを基に RDF テンプレートに記述

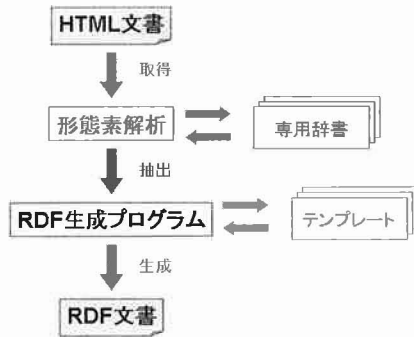


図2：提案手法

3.1. 形態素解析による要素抽出（作業1, 2）

形態素解析とは、入力された文字列を単語ごとに分解し品詞などの文法情報や意味情報など付加する処理である。本研究では、形態素解析のソフトとして『茶筌』を使い、単語ごとに分解し、自作した専用辞書を利用して、意味情報を付加する。

この際、意味情報を表現するために<タグ>を利用する。HTMLやXMLと同様に、タグで囲むことにより、HTMLパーサでの処理を可能とする。

3.2. RDF テンプレートの利用（作業3）

形態素解析での処理作業を通じて、必要な単語要素をHTMLと同様なタグで囲んだことによって、HTMLパーサを使用した解析処理が可能となった。パーサの役割としては、タグごとに分解された文章のタグが、形態素解析によって付加されたタグであるかのチェックを行い、タグで挟まれている要素をメタデータの目的語の候補として選択することである。また、囲まれているタグの種類によって、その要素がどの述語の目的語であるかも定義可能となる。

本研究では、パーサによって切り出したタグ情報を、基本となるRDFテンプレートに挿入することで新たにRDFファイルを作成する。

4. 検証実験

4.1. 実験準備

構築したRDF自動作成システムについて、検証実験を行った。現在使われているRDFには、Blog上でのRSSやDublinCoreによる書籍情報の管理などがあるが、本研究では、料理情報の管理についてRDFの自動作成を行う。

まず、料理情報に関する形態素解析用の辞書を作成した。形態素解析により付加するタグの意味情報は以下の4つを用意した。

- ◆ <ryouri> ～ 料理名
- ◆ <sozai> ～ 料理に使用される素材名
- ◆ <fumi> ～ 料理の風味の種類
- ◆ <tani> ～ 材料単位

これらのタグとRDFテンプレートを利用して、料理管理用のRDFファイルを作成した。

4.2. 実験結果

実証実験に使用したWebページを図3に示す。



図3：処理前のWebページ

これから、形態素解析、パーサ、テンプレートへの埋め込みの各作業を通じて作成されたRDFファイルの一部を以下に示す。

```

<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:eg="C:/sample/test.html">
  <rdf:Description rdf:about="C:/sample/test.html">
    <eg:ryouri xml:lang="ja">お好み焼き</eg:ryouri>
    <eg:fumi xml:lang="ja">和風</eg:fumi>
    <eg:image><imgsrc="http://image.infoseek.co.jp/
      recipe/1/199611326.jpg"width="220"height="169"border=
        "0"alt="広島風
    </eg:image>
    <eg:cal>723kcal</eg:cal>
    <eg:sozai rdf:parseType="Collection">
      . . . . (中略) . . . .
    <rdf:Description xml:lang="ja" rdf:about="#キャベツ"/>
    <eg:tani xml:lang="ja">550g</eg:tani>
    <rdf:Description xml:lang="ja" rdf:about="#もやし"/>
    <eg:tani xml:lang="ja">160g</eg:tani>
    . . . . (中略) . . . .
  </eg:sozai>
</rdf:Description>
</rdf:RDF>
  
```

5. おわりに

本研究では、これまでの情報処理の流れに大きな変革をもたらすであろう、セマンティックWebの容易な導入のために、RDFの自動作成ということについて議論した。そして、現在あるWeb上の情報(HTML)から、形態素解析やXMLパーサなど既存の技術を利用して、RDFを自動作成可能なことを検証実験により示した。

参考文献

- [1] 南野朋之, 鈴木泰裕, 藤木稔明, 奥村学, “blogの自動収集と監視,” 人工知能学会誌, Vol.19, No.6, pp.511-520, 2004.
- [2] “MyRSS.jp,” <http://myrss.jp/>.